## DEMONSTRATION OF FORMULÆ FOR TRUE MEASUREMENT OF CORRELATION.

### By C. SPEARMAN.

It seems daily more evident, that one of the most important tasks awaiting psychologists is the accurate measurement of the 'correlations' (*i. e.*, the tendencies to concurrent variation) between psychical events, qualities, faculties or other characteristics. For this purpose a now well-known method of calculation has been evolved by Bravais, Galton and Pearson, which furnishes a numerical 'coefficient' measuring precisely the degree of proportionality between any two series of values; this coefficient is usually denoted by the symbol r.[1] Should on any occasion the correlation between the two series take some more complicated form than that of simple proportion, then r has to be supplemented by further terms to express such correlation completely; r remains, however, the principal term (with or without some unessential modification of outward shape).

Unfortunately, there is a considerable step between arriving at a coefficient of correlation and discovering the *true* coefficient. In the first place, the values immediately attainable by investigation are not those of the characteristics really investigated, but only those of measurements, and—in the case of psychology—for the most part very fallible measurements. Secondly, it is usually quite impossible to keep the investigation clear of many factors that do not properly belong to it. The actual effect of these two disturbances, the observational errors

---

[1] For method of calculation see this *Journal*, 1904, XV, pp. 77-8 (but for "median" substitute "average").

and the irrelevant factors, is not merely to diminish somewhat the accuracy of the calculation, but to render the apparent correlation (whether calculated or merely casually inspected) wholly untrustworthy. A large correlation may be obliterated; an illusive one may be conjured up where none exists really; it may even happen that a positive correlatiou is turned into an apparently negative one, or *vice versa.*

Now, two formulæ were given by me in this *Journal* some time ago,[1] whereby the effect of both these disturbances can, as I believe, be eliminated. The formulæ were, however, not accompanied by any *proofs.* So many other mathematical formulæ were given at the same time, that the formal demonstrations of them all would have made the article exceedingly cumbersome.[2] But since then I have repeatedly been asked for these proofs; some mathematicians have gone so far, as to doubt whether such proofs could possibly be valid. It therefore seems advisable to publish them.[3]

For convenience of demonstration, I will commence with the formula for eliminating irrelevant factors, although in application the other formula must be used first (to *all* the coefficients entering into the former formula).

I. *Proof of the formula for eliminating the effect of irrelevant factors.*[4]

Let X, Y, and Z denote the values of any three variable and correlated characteristics of objects of any particular class (for instance, their height, length and breadth). Let their average values be denoted by $a_x$, $a_y$ and $a_z$ respectively; let $a_x — X = x$, $a_y — Y = y$, and $a_z — Z = z$. Further, let $b_{xy}$ and $b_{xz}$ be such values, that $\Sigma \left[ x — (b_{xy} y + b_{xz} z) \right]^2$ is a minimum. Equating to o the differentials of this sum with respect to both $b_{xy}$ and $b_{xz}$, and solving these two equations for $b_{xy}$, we find

$$b_{xy} = \frac{\Sigma xy \cdot \Sigma z^2 — \Sigma xz \cdot \Sigma yz}{\Sigma y^2 \cdot \Sigma z^2 — (\Sigma yz)^2} = \frac{r_{xy} — r_{xz} \cdot r_{yz}}{1 — r_{yz}^2} \cdot \frac{\sqrt{\Sigma x^2}}{\sqrt{\Sigma y^2}}$$

where $r_{xy}$ denotes Pearson's correlational coefficient

$$\frac{\Sigma\, xy}{\sqrt{\Sigma\, x^2.\ \Sigma\, y^2}},$$

and $r_{xz}$ and $r_{yz}$ have similar meanings.

Analogously $\quad b_{yx} = \dfrac{r_{xy} - r_{xz} \cdot r_{yz}}{1 - r_{xz}^2} \cdot \dfrac{\sqrt{\Sigma\, y^2}}{\sqrt{\Sigma\, x^2}}$, so that

$$\sqrt{b_{xy} \cdot b_{yx}} = \frac{r_{xy} - r_{xz} \cdot r_{yz}}{\sqrt{(1 - r_{xz}^2)\ (1 - r_{yz}^2)}} \qquad (a)\ [1]$$

Let us next suppose our whole class of objects to be split up into groups, a group embracing all those objects for which Z has any constant value. Let us apply the considerations and notation of the preceding paragraph to the $k^{th}$ such group apart. The term $b_{xz_k}$ vanishes, since $z_k$ clearly $= 0$.

We get, then, $\Sigma(x - b_{xy_k} y)^2 = $ a minimum, from which

$$\frac{d}{d\, b_{xy_k}} \Sigma\, (x - b_{xy_k} y)^2 = 0, \text{ so that } b_{xy_k} = \frac{\Sigma\, xy_k}{\Sigma\, y_k^2}$$

As $b_{yx_k}$ has an analogous value,

$$\sqrt{b_{xy_k} \cdot b_{yx_k}} = \frac{\Sigma\, xy_k}{\sqrt{\Sigma\, x_k^2.\ \Sigma\, y_k^2}} = r_{xy_k} . \qquad (b)$$

Now, in general the value of $b_{xy_k}$ differs from that of $b_{xy}$; but in three special cases here interesting us it can be shown to coincide.

The first case occurs when the following assumptions are permissible:

$$b_{xy_0} = b_{xy_1} = b_{xy_2} = \ldots = b_{xy_k} = \ldots \, , \qquad (c)$$
$$b_{yx_0} = b_{yx_1} = b_{yx_2} = \ldots = b_{yx_k} = \ldots \, , \qquad (d)$$
$$a_{x_k} - a_{x_0} = e \cdot Z_k \text{ and} \qquad (e)$$
$$a_{y_k} - a_{y_0} = f \cdot Z_k \, , \qquad (f)$$

e and f being constants. It will be convenient to conceive all the objects as represented geometrically by positions having as ordinates X, Y and Z. The equations $x - b_{xy_k} y = 0$ evidently denote what may be called 'minimal lines,' $b_{xy_k}$ being determined for each value of k by the relation $\Sigma\, (x - b_{xy_k} y)^2$ $= $ a minimum ; and all such minima may be regarded as part

---

[1] This result was reached, in a somewhat different manner, by Yule (Proc. R. Soc., Vol. LX).

sums of a total minimum which, transforming from x and y to X and Y, may be written as

$$\sum_k \sum \left[ X - b_{xy_k}Y - (a_{x_k} - b_{xy_k} a_{y_k}) \right]^2,$$

or, shorter, as $\sum_k \sum (X - b_{xy_k}Y - E_k)^2$, or simply as $\sum m$.

On the other hand, $x - b_{xy}y - b_{xz}z = 0$ is a 'total minimal plane,' $b_{xy}$ and $b_{xz}$ being determined by the relation $\sum (x - b_{xy}y - b_{xz}z)^2 = $ a minimum ; and this minimum can be regarded as made up of part sums (not necessarily minima individually), one for each different value of Z, and therefore may be written as

$$\sum_k \sum \left[ X - b_{xy}Y - (b_{xz}Z^k + a_x - b_{xy}a_y - b_{xz}a_z) \right]^2,$$

or, shorter, as $\sum_k \sum (X - b_{xy}Y - E^k)^2$,

or simply as M. Evidently, the difference between $\sum m$ and M depends on that between $E_k$ and $E^k$. But the former is the value that makes every part sum of the form in question a minimum, as may readily be found by determining E from the equation $\dfrac{\delta}{\delta E}(X - bY - E)^2 = 0$. Hence, if $E^k$ differs from $E_k$ for any single value of k, the corresponding part sum of M becomes greater than that of $\sum m$. *A fortiori*, the total effect of all differences between $E^k$ and $E_k$ for all values of k must be to make M greater than $\sum m$. But in the present case it happens to be possible for all corresponding parts of $\sum m$ and M respectively to be identical ; for the 'minimal lines,' from which $\sum m$ derives, must—owing to the conditions (c), (e) and (f)—lie in one plane ; and there is no condition preventing this plane from coinciding with the 'total minimal plane,' from which M derives. Since M *can* wholly coincide with $\sum m$, it *must* do so, for otherwise it would not be a minimum, as required by hypothesis. Hence finally, $b_{xy} = b_{xy_k}$ and, taking into consideration the equations (a) and (b), we arrive at the desired result,

$$r_{xy_k} = \frac{r_{xy} - r_{xz} \cdot r_{yz}}{\sqrt{(1 - r_{xz}^2)(1 - r_{yz}^2)}} \qquad (g)$$

The second case is where equation (f) is assumed once more and also $b_{xy_0} = b_{xy_1} = b_{xy_2} = \ldots = b_{xy_k} = \ldots = 0$. (h) It is clear that all the minimal lines must again lie in one plane, and this can again be shown, by reasoning as above, to coincide necessarily with the 'total minimal plane.' Consequently $b_{xy} = b_{xy_k} = 0$. Hence equation (g) again holds good, either side of it now being equal to 0.

The third case is similar to the first, with the exception, that $b_{xy_k}$ is permitted to vary with k. Under these conditions the equation (g) will be found valid to the extent of giving the most probable mean value of $r_{xy_k}$.

In equation (g) we have reached at any rate the form of the formula for eliminating irrelevant factors. But now we have to see whether, also, the constituent terms of this equation correspond to the respective values concerned in the actual investigation of correlation. Let us therefore turn to actuality, say, to the correlation found in my previous paper between children's power of discriminating pitch and that of discriminating weight. As irrelevant factor let us take the children's varying age. Clearly, the actually observed correlational coefficients are derived from three variables, discrimination of pitch, discrimination of weight, and age; they correspond perfectly to $r_{xy}$, $r_{xz}$ and $r_{yz}$ on the right of equation (g). But the factor of age is obviously irrelevant and disturbing; suppose, for instance, that both discrimination of pitch and that of weight improved with age; then of course we should find a correlation between the two sorts of discrimination, for both would be best in the same children, namely, the oldest; such correlation is evidently beside the question. Our true course of investigation would have been to select for experiment children of exactly the same age; but this is precisely the way we arrived at $r_{xy_k}$: thus $r_{xy_k}$ is the desired *true* correlation between x and y.

It remains to be considered how far the actual observational material fulfils the special conditions under which alone the formula (g) has been shown to be exact. The present topic, that is, the elimination of irrelevant factors, comes under our above 'first case;' we need, therefore, the relations, (c), (d), (e) and (f). Of these the two latter are equivalent to requiring the correlations of the irrelevant term (here, age) with the two main terms (here, discrimination of pitch and of weight) to be linear. This limitation is less serious than might be supposed; for the inexactness of the corrective formula only becomes appreciable when the irrelevant correlations depart from the linear form very largely; and experience has shown such large deviations to be extremely rare. Anyhow, it can scarcely be a matter of surprise that irrelevant correlations become difficult to treat when non-linear, seeing that no quite satisfactory formulæ have yet been discovered even for the bare measurement (*i. e.* antecedently to all corrections) of the main correlation when non-linear.

Finally, the two other conditions (c) and (d), mean that the true correlation must not change appreciably for the different values of the irrelevant term. Now, such changes may be taken

as, in general, of a smaller order of magnitude than the altera-
tions to be eliminated by the corrective formula, those produced
by *mixing* different values of the irrelevant term.  To return
to our instance, there is good reason to believe that most cor-
relations are very similar for children of 9 years as for those of
12, although the gravest disturbance will often occur if 9 and
12 years be thrown together into one and the same correlation.
Conditions (c) and (d) may therefore be considered sufficiently
satisfied whenever the irrelevant influence to be eliminated is
of moderate amount.  But if, instead of confining our experi-
ments to the ages of 9–12, we had included those down to, say,
5, then the true correlation for 5 years would probably have
had quite a large discrepancy from that for 12.  In such case
one could at most expect any general 'true' correlation to sig-
nify the true *mean* correlation; and this, as we have seen, is
the value actually given by our formula.

2. *Proof of the formula for eliminating the effect of inaccurate
observation.*

We will assume any two correlated series of values, X and Y
to have each been measured twice independently, and to have
yielded the series of measurements $x_1$, $x_2$, $y_1$ and $y_2$.  The co-
efficients $r_{x_1y_1}$, $r_{x_1y_2}$, $r_{x_2y_1}$, $r_{x_2y_2}$, $r_{x_1x_2}$ and $r_{y_1y_2}$ can, of course, be
reckoned directly.  We require a formula to reckon $r_{XY}$.

Let us first consider the correlations between $x_1$, $x_2$ and X,
and see how the coefficient between $x_1$ and $x_2$ becomes modi-
fied when a separate calculation is made for each group of ob-
jects for which X is constant, say, $= X_k$.  We may fairly as-
sume the average of all the measurements $x_{2k}$ (or $x_{1k}$) to co-
incide with $X_k$, or at any rate to vary proportionally thereto
for the different values of k; hereby the condition (f) is satis-
fied.  Further, when we consider any $k^{th}$ group quite apart,
since the fluctuations in the two series of measurements $x_{1k}$
and $x_{2k}$ of the same value $X_k$ are by hypothesis independent
of one another, $b_{x_1x_2k}$ (or $b_{x_2x_1k}$) always $= 0$, thus satisfying
the condition (h).  Consequently, we have here our 'second
case' and

$$r_{x_1x_2k} = \frac{r_{x_1x_2} - r_{x_1X} \cdot r_{x_2X}}{\sqrt{(1 - r^2_{x_1X})(1 - r^2_{x_2X})}} = 0, \quad \text{(i)}$$

where X is fixed at $X_k$ on the left side of the equation, but
remains variable on the right.  Therefore, since r cannot be
infinitely great,

$$r_{x_1x_2} = r_{x_1X} \cdot r_{x_2X}, \text{ and analogously} \quad \text{(j)}$$

$$r_{y_1y_2} = r_{y_1Y} \cdot r_{y_2Y} \quad \text{(k)}$$

Next, let us consider the correlations between $x_1$, X and $y_1$. If we again make a separate calculation for each group of values for which X is constant, the condition (f) is satisfied just as before. Further, in the calculation for any $k^{th}$ group considered quite apart, $X_k$, being a constant, is independent of the fluctuations in the measurements $y_{1_k}$, so that $b_{y_1 x_k} = 0$ and condition (h) is satisfied again. Hence our 'second case' occurs once more and

$$r_{x_1 y_1 k} = \frac{r_{x_1 y_1} - r_{x_1 X} \cdot r_{y_1 X}}{\sqrt{(1 - r^2_{x_1 X})(1 - r^2_{y_1 X})}} = 0. \qquad (1)$$

From this evidently

$$r_{y_1 X} = \frac{r_{x_1 y_1}}{r_{x_1 X}} \text{ and, analogously, } = \frac{r_{x_2 y_1}}{r_{x_2 X}}. \qquad (m)$$

Likewise $\qquad r_{y_2 X} = \frac{r_{x_1 y_2}}{r_{x_1 X}} = \frac{r_{x_2 y_2}}{r_{x_2 X}}, \qquad (n)$

$$r_{x_1 Y} = \frac{r_{x_1 y_1}}{r_{y_1 Y}} = \frac{r_{x_1 y_2}}{r_{y_2 Y}}, \qquad (o)$$

$$\text{and} \quad r_{x_2 Y} = \frac{r_{x_2 y_1}}{r_{y_1 Y}} = \frac{r_{x_2 y_2}}{r_{y_2 Y}}. \qquad (p)$$

Finally, let us take the correlations between $x_1$, X and Y. By reasoning as before, we get

$$r_{x_1 Y_k} = \frac{r_{x_1 Y} - r_{x_1 X} \cdot r_{X Y}}{\sqrt{(1 - r^2_{x_1 X})(1 - r^2_{X Y})}} = 0,$$

from which it is evident that

$$r_{X Y} = \frac{r_{x_1 Y}}{r_{x_1 X}} \text{ and analogously,}$$

$$= \frac{r_{x_2 Y}}{r_{x_2 X}} = \frac{r_{y_1 X}}{r_{y_1 Y}} = \frac{r_{y_2 X}}{r_{y_2 Y}}.$$

By multiplying together the four preceding equations to $r_{X Y}$, we have

$$r_{X Y}^4 = \frac{r_{x_1 Y} \cdot r_{x_2 Y} \cdot r_{y_1 X} \cdot r_{y_2 X}.}{r_{x_1 X} \cdot r_{x_2 X} \cdot r_{y_1 Y} \cdot r_{y_2 Y}}$$

Substituting on the right of the above equation from (j), (k), (m), (n), (o), and (p) and taking the real positive root, we find at last

$$r_{X Y} = \frac{G\left(r_{x_1 y_1}, r_{x_1 y_2}, r_{x_2 y_1}, r_{x_2 y_2}\right)}{G\left(r_{x_1 x_2}, r_{y_1 y_2}\right)}, \qquad (q)$$

where G denotes the geometrical mean.

In practice it will usually be allowable to assume that the two series of measurements of the same series of things have been conducted with equal accuracy. Then $r_{x_1 x} = r_{x_2 x}$ and $r_{y_1 y} = r_{y_2 y}$, so that equation (q) becomes

$$r_{x\,y} = \frac{r_{x\,y}\,(\,= r_{x_1 y_1} = r_{x_1 y_2} = r_{x_2 y_1} = r_{x_2 y_2}\,)}{G\,(\,r_{x_1 x_2}\,,\,r_{y_1 y_2}\,)}.$$

The discrepancies that will occur between the four *actual* values of $r_{x\,y}$ must then be attributed to mere chance, and must be met, as usual, by taking an average. Thus we get, on the assumption of the equally accurate series of measurements,

$$r_{x\,y} = \frac{\dfrac{r_{x_1 y_1} + r_{x_1 y_2} + r_{x_2 y_1} + r_{x_2 y_2}}{4}}{\sqrt{r_{x_1 x_2} \cdot r_{y_1 y_2}}}.$$

The above proof of the formula for eliminating the effect of observational errors is, as we have seen, exact and perfectly general. It holds good whatever may be the distribution of values, or the size or distribution of the observational errors, in the series concerned and whatever may be the correlation's form. Any discrepancies arise solely from the practical necessity of applying the formula, not to the whole series of values considered, but only to 'random samples' of such series. By sufficiently extending the experiments, the chance of discrepancy may be reduced as much as desired. The formula for irrelevant factors is equally general, except for the two limitations explained above.

Both formulæ concern themselves, however, solely with r, that is, with $\dfrac{\Sigma xy}{\sqrt{\Sigma x^2 \cdot \Sigma y^2}}$. But, as mentioned before, when the relation between the two characteristics investigated assumes some special form, instead of the normal 'linear' form of simple proportion, then this special form finds no expression whatever in r taken alone. To express this form analytically, other additional terms are required. The exact nature of these additional terms (as well as the outward shape of r itself) varies somewhat according to the method of calculation adopted. But there is, probably, no mathematical difficulty in devising modifications of our corrective formulæ to suit any such terms.

It should be observed, that in many cases the non-linear form is more apparent than real. Generally speaking, a mere tendency of two characteristics to vary concurrently must be taken, it seems to me, as the effect of some particular underlying strict law (or laws) partly neutralized by a multitude of 'casual' disturbing influences. The quantity of a correlation is neither more nor less than the relative influence of the un-

derlying law in question as compared with the total of all the influences in play.   Now, it may easily happen, that the underlying law is one of simple proportionality but the disturbing influences become greater when the correlated characteristics are larger (or smaller, as the case may be).   Then the underlying simple proportionality will not appear on the surface; the correlation will seem non-linear.   Under such circumstances, r cannot, it is true, express these variations in the quantity of correlation; it continues, however, to express completely the *mean* quantity of correlation.

In the majority of the remaining cases of non-linearity, the latter is merely due to a wrong choice of the correlated terms. For instance, the correlation between the length of the skull and the weight of the brain must, obviously, be very far from linear.   But linearity is at once restored (supposing all the skulls to belong to one type) if we change the second term from the brain's weight to the cube root of the weight.

To conclude, even when the underlying law itself really has a special non-linear form, although r by itself reveals nothing of this form, it nevertheless still gives (except in a few extreme and readily noticeable cases) a fairly approximate measure of the correlation's quantity.[1]

---

[1] Several writers, who have made otherwise valuable contributions to the subject of correlation, but have been too exclusively guided by the purely mathematical point of view, appear to have wholly overlooked this fundamental distinction between the form and the quantity of a correlation.